

Regularization

L2 regularization:

Logistic regression:

$$\min_{w, b} J(w, b)$$

loss function.

λ = regularization param.

$$\rightarrow J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \|w\|_2^2$$

$$\text{L2 regularization: } \|w\|_2^2 = \sum_{j=1}^{n_x} w_j^2 = \underline{w^T w^T}$$

Result will be a scalar.

$$\text{L1 regularization: } \frac{\lambda}{2m} \sum_{j=1}^{n_x} |w_j| = \frac{\lambda}{2m} \|w\|_1.$$

Neural Network:

$$J(w^{[0]}, b^{[0]}, \dots, w^{[L]}, b^{[L]}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{[l]}\|_F^2$$

Frobenius norm:

$$\|w^{[l]}\|_F^2 = \sum_{i=1}^{n_l} \sum_{j=1}^{n_{l-1}} (w_{ij}^{[l]})^2$$

$$\rightarrow dw^{[l]} = (\text{from backprop}) + \frac{\lambda}{m} w^{[l]}$$

$$w^{[l]} = w^{[l]} - \alpha dw^{[l]}$$

$$\Rightarrow w^{[l]} = w^{[l]} - \alpha \left[(\text{from backprop}) + \frac{\lambda}{m} w^{[l]} \right]$$

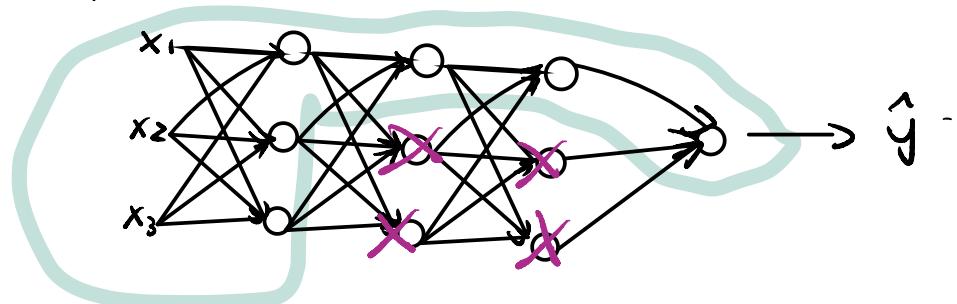
$$= w^{[l]} - \alpha (\text{from backprop}) - \frac{\alpha \lambda}{m} w^{[l]}$$

$$= \underbrace{\left(1 - \frac{\alpha \lambda}{m}\right)}_{\text{Normal weight decay}} w^{[l]} - \alpha (\text{from backprop}).$$

Normally less than 1. \rightarrow "weight decay"

Why regularization reduces overfitting

Explanation 1:



$$J(w^{[l]}, b^{[l]}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{[l]}\|_F^2$$

Minimizing J needs $\underline{w^{[l]} \approx 0}$

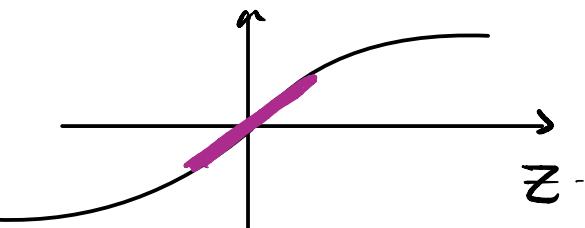
$w^{[l]} \approx 0 \rightarrow$ Some neuron (hidden units)' effects are diminished, if not shut down

Explanation 2:

e.g. tanh: $g(z) = \tanh(z)$.

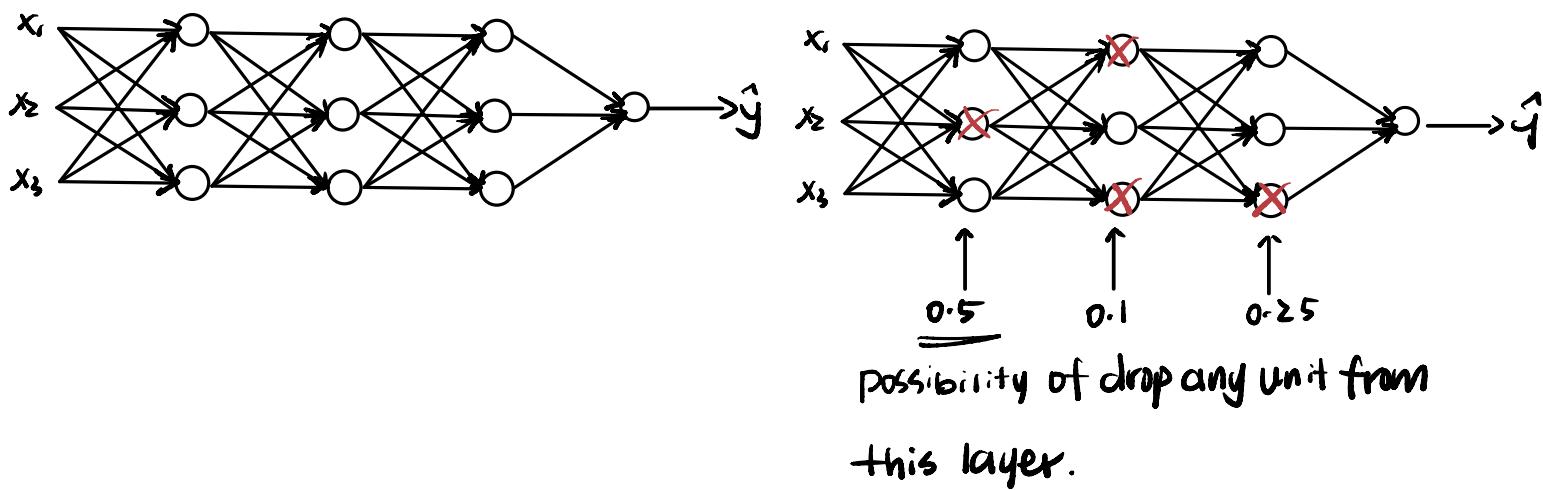
$$z^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]}$$

As shown above, if $w^{[l]}$, narrowing z into the small range shown in purple line, which makes each layer perform like a linear function. Thus reduce variance. (overfitting).



Dropout regularization.

Intuition: randomly shutting down a small portion of hidden units in each iteration.



Implementation ("Inverted dropout").

Illustrate with layer $l = 3$, keep-prob = 0.8

$\rightarrow d_3 = \text{np.random.rand}(a_3.\text{shape}[0], a_3.\text{shape}[1]) < \text{keep_prob}$

$a_3 = \text{np.multiply}(a_3, d_3)$. drop out.

$a_3 / \equiv \text{keep_prob}$.

Keep a_3 's original magnitude.

(Since 20% are shut off, the remaining need to $\times 0.8$).

$$\rightarrow (-0.2)X = 0.8X$$

$$0.8X / 0.8 = \underline{\underline{X}}$$

* NO dropout when making predictions.

Potentially you can dropout multiple times and calculate their mean, using that as predictions for test set.

However, this is very expensive computation which normally produce similar results.

why does drop-out work?

→ can't rely on any one feature. so have to spread out weights

resulting in a more generalized model, reducing variance.

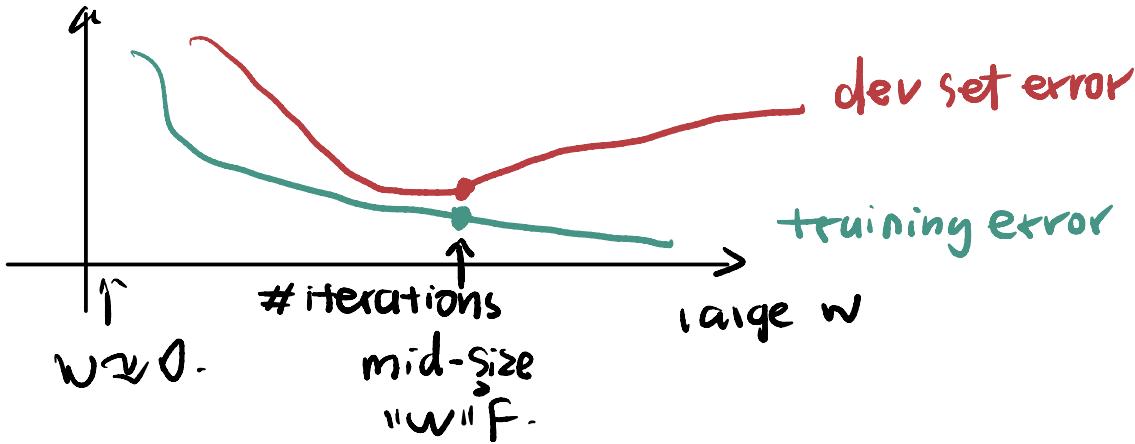
More regularization methods.

Data Augmentation:

e.g. Computer Vision:

through distortion, rotation, reversion, transform pictures in training set to increase training size.

Early Stopping



* Orthogonalization: Achieve two goals below one after another.

→ - Optimize cost function J
- Gradient Descent ...

→ - Not overfit
- Regularization ...

early stop can be sometimes problematic because it obeys orthogonalization. Thus, it makes trade-off btw bias and variance.